

工學碩士學位論文

A Study on Learning a Text Filtering System Using  
Boosting Algorithm.

2000年 2月

大學校 大學院

工學科

韓 尚 潤

A Study on Learning a Text Filtering System Using  
Boosting Algorithm.

指導教授 金 榮 澤

論文 工學碩士 學位論文 提出

1999年 10月

大學校 大學院

工學科

韓 尚 潤

韓尚潤 工學碩士 學位論文 認准

1999年 12月

委 員 長 \_\_\_\_\_ 印

副委員長 \_\_\_\_\_ 印

委 員 \_\_\_\_\_ 印

(text filtering)

가

가

가

가

.

(AdaBoost)

가

가

가

가

1

-1

가

가

가

(continuous Poisson distribution)

TREC (Text Retrieval Conference)

TREC-8 adhoc collection

Financial Times 1992 1994

가

가

가

TREC 가

:

,

,

,

,

1	.....	4
1.1	.....	4
1.2	.....	6
1.3	.....	7
2	.....	8
2.1	.....	8
2.2	.....	10
2.3	.....	14
3	.....	15
4	.....	18
4.1	.....	18
4.2	가 .....	21
4.3	.....	25
5	가 .....	27
5.1	.....	27
5.2	가 .....	31
6	.....	42
	.....	44

1	AdaBoost.M1	.....	11
2		.....	19
3	가	.....	22
4		.....	25
5	Reuter	.....	27
6	TREC_8 ad hoc collection	.....	30
7	real- AdaBoost    i- AdaBoost    LF1	.....	36
8		.....	38
9	TREC- 8    가        i- AdaBoost    LF1	.....	40
10		.....	41

1 AdaBoost- tf real- AdaBoost	32
2 AdaBoost- tf	33
3 real- AdaBoost	33
4 AdaBoost- tf real- AdaBoost	34

# 1

## 1.1

가

.

.

(1 -1)

가

가

가

가 , , ,

가

가

가

가 . ,

가

가

가

가 가 , 가  
가 가 1992  
TREC (Text REtrieval Conference)

가

가 가

가

$k$

[lew is92, lew is94,

yang94, yang95, hull96, lewis96].

(boosting)

[schapire98b]

가

(decision- committee model)

(bagging)

가

가

(bias), (variance)

[schapire98a]

[bauer99].

가

(overfitting)

가

[schapire98b]



## 1.2

가

가

가

가

가

TREC (Text REtrieval Conference)

, 가 . TREC 1992

8

, , 7

8

가

50

가

가

가

### 1.3

2

PAC(Probably Approximately correct)

,

. 3

가

. 4

, 5

,

.

## 2

### 2.1

가  
 , (strong learning model) (weak  
 learning model) . Valiant 1984 가  
 PAC (Probably Approximately Correct)  
 [neuralnet].  $x$ 가  $x = (a_1, \dots, a_n)$   
 가  $X$  . 가  
 $C, c(x) \in \{0, 1\}$   $x$ 가 .  
 $D$   $X$   
 $x$   $c(x) = h(x)$ 가  $h(x)$  .  
 가 0 가  $h$   
 가 . ,  
 가 , 가  
 가 0 가 . , 가  
 $\epsilon$  , 가  
 가  $\delta$ 가 가 가  
 PAC . PAC  
 .  
 가 [neuralnet].

가 1/2 가  
 " 가"  
 Schapire 1990 [schapire90]  
 가  
 . n 가  
 가 . 1 가  
 n 가 . 가  
 가 n 가  
 가 가 , 가  
 가 가 가  
 가 ε 3ε<sup>2-2ε<sup>3</sup></sup>  
 [schapire90]

## 2.2

가

가

[freund96a].

가

가

가

1

AdaBoost.M1

[freund96a].

가 (uniform) 가

가

$1/m$

가

가

가  $1/2$

가  $1/2$

가

가

가

가

가

가

, 가

가

가

$\log(1/\beta)$

가

$m$  가  $S$ ,  
 $I$ ,  
 $T$   
1  $s' = S$  (  $\frac{1}{m}$  )  
2 For  $i = 1$  to  $T$  {  
3  $C_i = I(S')$   
4  $\epsilon_i = \frac{1}{m} \sum_{x_j \in S': C_i(x_j) \neq y_j} \text{weight}(x)$  (  $\frac{1}{m}$  )  
5  $\epsilon_i > \frac{1}{2}$   $S$   $S'$   
 $\frac{1}{m}$  3  
6  $\beta_i = \frac{\epsilon_i}{1 - \epsilon_i}$   
7  $x_j \in S'$  ,  $C_i(x_j) = y_j$   $\text{weight}(x_j) \cdot \beta_i$   
 $\frac{1}{m}$   
8  $S'$   $\frac{1}{m}$  가  
9 }  
10  $C^*(x) = \arg \max_{y \in Y} \sum_{i: C_i(x) = y} \log\left(\frac{1}{\beta_i}\right)$  :  $C^*$

## 1 AdaBoost.M1

$n$  가 (가 )  
 $D_n$ ,  $n$  가  $\epsilon_n$  ( $\epsilon_n < (1/2)$ ),  $\gamma_n = (1/2) - \epsilon_n$

$$\frac{|\{i: C^*(x_i) \neq y_i\}|}{m} \leq \prod_{n=1}^T \sqrt{1 - 4\gamma_n^2} \leq \exp(-2\sum_{n=1}^T \gamma_n^2)$$

가 가 가

Schapire

(margin)

[schapire98a].

가 가 가 가 가

가

Breiman

가

가

가

[breiman97].

가

가

C4.5 [quinlan96], perceptron [freund99], Naive Bayes

[bauer99]

,

[freund96b],

[abney99],

[schapire98b, schapire99b]

가 가  $\alpha$

[freund97],

가 -1 1

[schapire99a].

가

4

가

가

,  
[friedman]  
[breiman96b, schapire98a]. 가  
가  
[freund96]. , 가 가  
가 .



### 2.3

(bagging) , 가

가 . 가

가 가

$n$  가 .

가 가

[breiman96a]. 가

(instable) [breiman96a].

가 [breiman96b].

가 [bauer99].

, (MultiBoosting)

[webb]. 가

### 3

가 가  
가

$$(1 \quad -1)$$

$d$  가

$y$

$$\langle (d_1, y_1), \dots, (d_m, y_m) \rangle$$

$$(d_i \in X, y_i \in \{-1, 1\})$$

$$C: X \rightarrow \{-1, 1\} \quad C$$

가 가

$$\epsilon = 1/m \sum_{i: C(d_i) \neq y_i} 1$$

TREC (Text REtrieval

Conference) 7

가

( )

가

가

가 .

가 1000 .

가 .

가 가 .

가 가 .

ε .

가 .

가 .

(precision) 가 . (recall)

	$R_+$	$N_+$
	$R_-$	$N_-$

$R_+$

,  $R_-$

,  $N_+$  ,  $N_-$

$r$

$$r = \frac{R_+}{R_+ + R_-} \quad p \quad p = \frac{R_+}{R_+ + N_+}$$

(break-even point)

가

가 .

[schapire98b].

가

가

$F_1$

$$F_1 = \frac{2rp}{r+p}$$

가

가

가

가가

가

TREC

$$aR_+ + bN_+ + cR_- + dN_-$$

$a, b, c, d$   $R_+, N_+, R_-, N_-$

가 .

TREC

$$LF1 = 3R_+ - 2N_+ \quad LF2 = 3R_+ - N_+ \text{ 가 .}$$

가

가

LF1



$D_{s+1}(i)$

가

$$\begin{aligned}
 & : N \quad \quad \quad : \\
 & \langle (d_1, y_1), \dots, (d_N, y_N) \rangle \quad (y_i \in \{-1, 1\}) \\
 & \quad \quad \quad T \quad \quad \quad \text{가} \quad \quad \quad ) \\
 & : \quad \quad \quad \text{가} \quad \quad \quad D \quad \quad \quad D_1(i) = \frac{1}{N} \\
 & \quad \quad \quad s = 1 \text{ to } T \\
 & 1 \quad \quad \quad D_s \quad \quad \quad \text{가} \\
 & \quad \quad \quad \text{가} \quad \quad \quad h_s \quad \quad \quad . \\
 & 2 \quad \quad \quad h_s \quad \quad \quad \epsilon_s = \sum_{i: h_s(d_i) \neq y_i} D_s(i) \quad \quad \quad . \\
 & 3 \quad \quad \quad \alpha_s = \frac{1}{2} \ln\left(\frac{1 - \epsilon_s}{\epsilon_s}\right) \\
 & 4 \quad \quad \quad . \\
 & D_{s+1}(i) = \frac{D_s(i) \exp(-\alpha_s y_i h_s(d_i))}{Z_s} = \frac{D_s(i)}{Z_s} \cdot \begin{cases} e^{-\alpha_s} & \text{if } h_s(d_i) = y_i \\ e^{\alpha_s} & \text{if } h_s(d_i) \neq y_i \end{cases} \\
 & \quad \quad \quad Z_s \quad \quad \quad . \\
 & \quad \quad \quad \text{가} \quad \quad \quad . \\
 & h_{fin}(d) = \text{sign}\left(\sum_{i=1}^T \alpha_s h_s(d)\right)
 \end{aligned}$$

2

$$\begin{aligned}
 & \text{가} \quad \quad \quad ( \quad \quad \quad ) \text{가} \quad \quad \quad \text{가} \\
 & \quad \quad \quad 1 \quad \quad \quad -1 \quad \quad \quad . \quad \quad \quad \text{가} \\
 & \quad \quad \quad \text{가} \quad \quad \quad \text{가} \\
 & \quad \quad \quad . \quad \quad \quad \text{가} \\
 & \quad \quad \quad , \quad \quad \quad \epsilon = \sum_{i: h_s(d_i) \neq y_i} D_s(i)
 \end{aligned}$$

$\epsilon$  가  $h_s$  가  
 가 1 , -1  
 가 (0.9) 가  
 가  
 $|\epsilon - (1 - \epsilon)|$  가 가  
 가 가

$$h_{fin}(d) = \text{sign}(\sum_{i=1}^T \alpha_i h_s(d))$$

가 가 가  $\alpha_s$   
 가  $h_s$  가 가 1/2 가  
 가  $h_s$  가  $\alpha_s$ 가 0 가  $h_s$  가  
 가 0 가  $\alpha_s$   
 가 , 가 1 가  
 $\alpha_s$ 가

[schapire98b]

Rocchio

가

Rocchio

## 4.2 가

Schapire Singer [schapire99a] 가  
 가  $h: X \rightarrow R$   
 $1 - 1$   
 가  
 $1, -1$   
 가  
 $|h|$   
 가  
 $1 - 1$  가  
 가

가  
 $f(x) = \sum_{s=1}^T \alpha_s h_s(x), \pi$   
 $I(\pi) \pi$  가  $1, -1$  가  
 $H(x) = \text{sign}(F(x))$

: 가  $H$

$\frac{1}{m} |(i: H(x_i) \neq y_i)| \leq \sum_{s=1}^T \alpha_s Z_s$  [schapire99a]  
 $\alpha_s$  가  
 가  $\alpha_s Z_s$   
 가  $\alpha_s$   
 가  
 $h_s \alpha_s$



$$\begin{aligned}
 & : \quad \langle (x_1, y_1), \dots, (x_m, y_m) \rangle \quad (x_i \in X, y_i \in \{-1, 1\}) \\
 & \quad \quad \quad T \\
 & : \quad \quad \quad \text{가} \quad D_{1(i)} = 1/m. \\
 \text{For } s = 1 \text{ to } T \\
 & \quad 1 \quad D_s \quad \quad \quad \text{가} \quad h_s: X \rightarrow R \\
 & \quad 2 \quad \alpha_s \in R \\
 & \quad 3 \quad \quad \quad : \\
 & \quad D_{s+1}(i) = \frac{D_s(i)}{Z_s} \exp(-\alpha_s y_i h_s(x_i)) \quad Z_s \\
 & \quad \quad \quad \cdot \\
 & \quad Z_s = \sum_i D_s(x_i) \exp(-\alpha_s h_s(x_i) y_i) \\
 & \quad \quad \quad \text{가} \quad \cdot \\
 & \quad H(x) = \text{sign}(\sum_{s=1}^T \alpha_s h_s(x)) \\
 & \quad \quad \quad 3 \quad \quad \quad \text{가}
 \end{aligned}$$

$$\begin{aligned}
 & \quad \quad \quad \alpha_s \quad 1 \quad \quad \quad \text{가} \\
 & \quad \cdot \quad \quad \quad , \quad Z_s = \sum_i D_{s(i)} \exp(-y_i h_s(x_i)) \quad h_s \\
 & \quad \cdot \quad \quad \quad \text{가} \quad \quad \quad x_1, \\
 & \quad \quad \quad X_0 \quad \quad \quad , \quad x \in X, j \in \{0, 1\}, \quad c_j = h(x) \\
 & \quad \cdot \quad \quad \quad \text{가} \quad \quad \quad \text{가} \quad \quad \quad \text{가} \\
 & \quad \text{가} \\
 & \quad w_b^j = \sum_{i: x_i \in X_j \wedge y_i = b} D(i) \quad \cdot \quad b \quad \text{가}
 \end{aligned}$$

+ - . Z

$$Z = \sum_j \sum_{i: x_i \in X_j} D(i) \exp(-y_i c_j) = \sum_j (w_+^j e^{-c_j} + w_-^j e^{c_j})$$

Z c\_j 가 0 c\_j

$$c_j = \frac{1}{2} \ln\left(\frac{w_+^j}{w_-^j}\right) . Z$$

Z = 2 \sum\_j \sqrt{(w\_+^j w\_-^j)} 가 h . c\_j

가

가 . 가 가 가 .  
c\_1 가 , 가 가 .

|c\_j| .

w\_+^j, w\_-^j 0 0 가 Z

c\_j 가 . |c\_j| 가 가 .  
D 가 .

w\_+^j w\_-^j E Z

c\_j .

$$Z = \sum_j w_+^j \sqrt{\left(\frac{w_-^j + E}{w_+^j + E}\right)} + w_-^j \sqrt{\left(\frac{w_+^j + E}{w_-^j + E}\right)}$$

$$c_j = \frac{1}{2} \ln\left(\frac{w_+^j + E}{w_-^j + E}\right)$$

$$|c_j| = \frac{1}{2} \ln\left(\frac{1}{E}\right)$$

. 4.1 가

4

가

가

가

가

가

### 4.3

가 . 가 가 .  
 가 가  
 가 ?

가

가 .

:  
 $T,$   
 $U$   
 1.  $T$  가 .  
 2.  $D$  (continuous Poisson distribution)  
 $1$   
 3.  $1, 2$   $U$  .  
 4. 가 가 가 가

$D_1(i)$  가

$1/m$

$D$

가

가

[webb] 가

가

$$poisson() = - \log \left( \frac{random(1 \cdots 999)}{1000} \right)$$

$random(min \cdots max)$     min    max

가

2

가

TREC

가

가

## 5 가

### 5.1

Reuter TREC\_8 adhoc  
 Reuter 10 가  
 8762 , 3009  
 8754 ( ) TF-IDF  
 4.1  
 4.2 가

acq	1483	7279	640	2369
corn	133	8629	36	2973
crude	334	8428	156	2835
earn	2706	6056	1034	1966
grain	370	8392	125	2884
interest	275	8487	97	2912
money- fx	428	8334	3009	131
ship	180	8582	80	2929
trade	303	8459	102	2907
wheat	185	8577	54	2955

### 5 Reuter

TREC-8 ad hoc collection , 1992 Financial Times  
 1993 1994  
 300 stop stemming  
 TF-IDF 300 351 400 50  
 가

가 , . ,

351	160	11	149	446	17	429
352	499	55	444	1109	182	927
353	152	15	137	267	24	243
354	212	21	191	319	41	278
355	83	1	82	208	1	207
356	255	8	247	484	7	477
357	143	24	119	229	40	189
358	53	0	53	106	0	106
359	186	4	182	500	17	483
360	66	2	64	152	12	140
361	181	0	181	412	1	411
362	61	0	61	129	3	126
363	164	2	162	506	3	503
364	113	1	112	210	2	208
365	142	8	134	205	3	202
366	143	4	139	465	15	450
367	149	12	137	263	25	238
368	87	87	0	262	5	257
369	104	1	103	232	0	232
370	107	15	92	232	6	226
371	125	1	124	407	1	406
372	99	3	96	247	12	235
373	144	2	142	359	10	349
374	97	19	78	264	53	211
375	106	9	97	198	4	194
376	363	13	350	591	15	576
377	136	5	131	270	10	260
378	467	51	416	813	36	777
379	118	0	118	292	0	292
380	104	0	104	273	1	272
381	136	5	131	328	2	326
382	70	2	68	170	4	166
383	114	15	99	358	65	293
384	57	1	56	131	3	128
385	125	9	116	213	25	188
386	107	1	106	235	3	232
387	55	4	51	220	9	211
388	96	3	93	170	12	158
389	237	41	196	458	98	360
390	151	2	149	332	49	283
391	221	55	166	673	106	567



392	189	27	162	322	21	301
393	244	4	240	305	5	300
394	151	0	151	396	5	391
395	186	36	150	341	51	290
396	101	6	95	283	3	280
397	123	2	121	325	5	320
398	145	10	135	191	7	184
399	81	4	77	143	6	137
400	112	34	78	149	16	133

## 6 TREC\_8 ad hoc collection

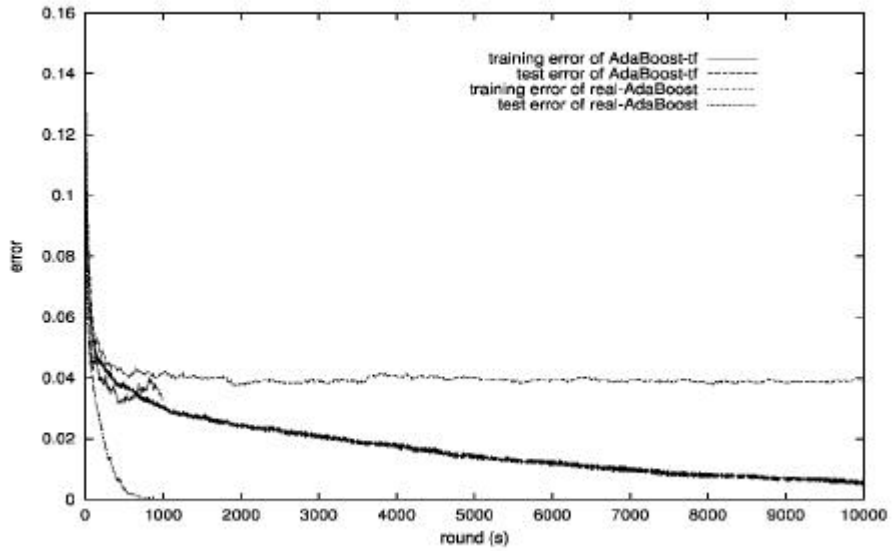
가

가

가

## 5.2 가

가  
 가 , ,  
 Reuter 1 acq  
 AdaBoost- tf (4.1 ) real- AdaBoost (4.2 )  
 .  
 AdaBoost- tf 10,000 가 ,  
 real- AdaBoost 1,000 가 .  
 real- AdaBoost가  
 . 가  
 AdaBoost- tf 4,000 real- AdaBoost 430  
 .  
 real- AdaBoost가 10  
 . 가  
 0  
 .  
 가 가



1 AdaBoost- tf real- AdaBoost

AdaBoost- tf 4.1

, real- AdaBoost 4.2

가

2 3 AdaBoost- tf real- AdaBoost

4

AdaBoost- tf 10,000 , real- AdaBoost 1,000

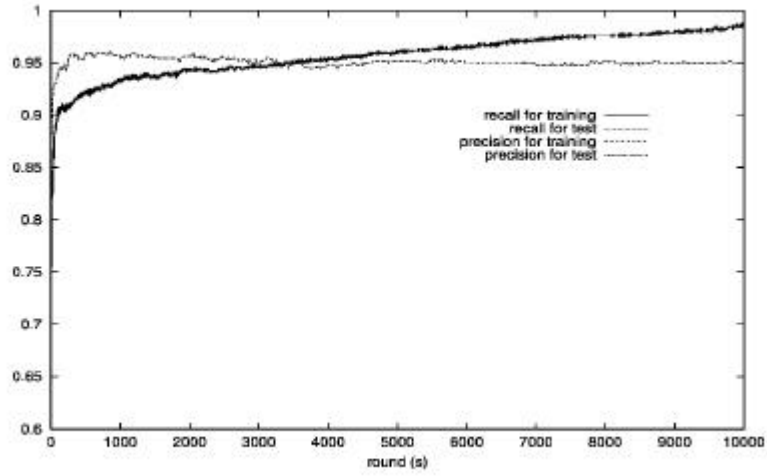
3

real- AdaBoost , 1,000

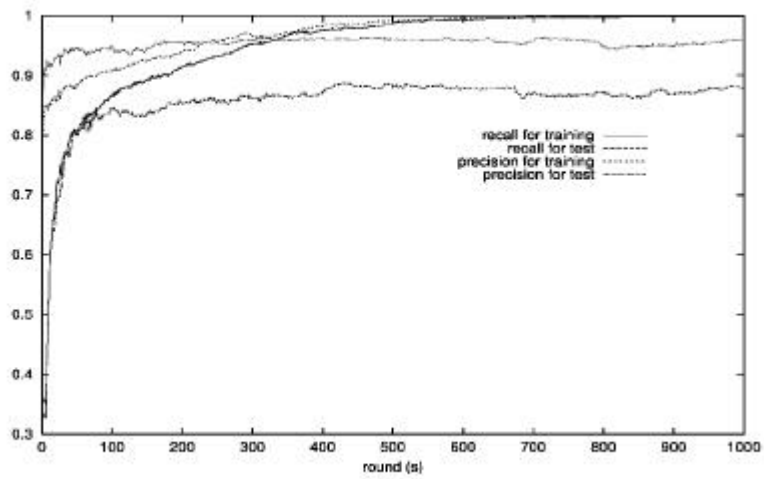
1 , 88%,

96% . 4

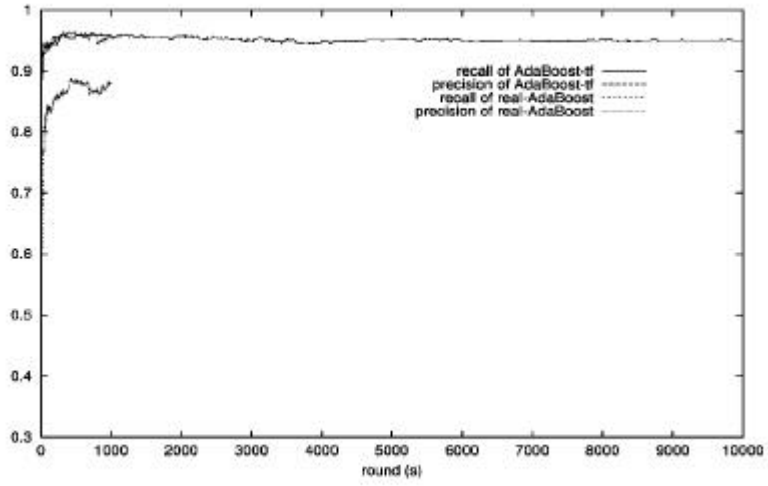
real- AdaBoost가 AdaBoost- tf



2 AdaBoost- tf



3 real- AdaBoost



4 AdaBoost- tf real- AdaBoost

1 4 real- AdaBoost가 AdaBoost- tf  
 가  
 real- AdaBoost 가  
 AdaBoost- tf 가  
 AdaBoost- tf 가 1,  
 - 1 가 가 가  
 $\alpha$  , real- AdaBoost 가  
 AdaBoost- tf 가  
 $|\alpha^* h_s|$  가 real- AdaBoost 가  
 가

TREC- 8

5 가

real- AdaBoost i- AdaBoost (4.3 )  
 real- AdaBoost 100 가 가 가  
 가 LF1 LF1 가 LF1  
 50 가 LF1  
 i- AdaBoost  
 가 LF1  
 30, 50, 70, 100  
 5 LF1 50  
 LF1  
 real- AdaBoost 7 8  
 30, 50, 70, 100

	real- AdaBoost		i- AdaBoost	
		LF1		LF1
351	25	63	62	12
352	21	20	25	63
353	9	13	21	20
354	16	0	68	14
355	9	1	16	0
356	3	34	9	1
357	1	0	3	34
358	1	0	1	0
359	2	1	1	0
360	1	0	2	1
361	1	0	1	0
362	1	0	1	0
363	1	0	1	0
364	1	0	1	0
365	1	0	1	0
366	7	6	107	3

	real- AdaBoost		i- AdaBoost	
		LF1		LF1
367	1	0	7	6
368	17	- 2	1	0
370	1	0	1	0
371	17	0	17	0
372	2	6	2	6
373	1	0	1	0
374	1	0	202	33
375	1	0	1	0
376	1	0	1	0
377	1	0	1	0
378	1	0	1	0
379	1	0	1	0
380	1	0	1	0
381	1	0	1	0
382	3	0	3	0
383	1	84	1	84
384	1	0	1	0
385	1	0	104	22
386	18	0	18	0
387	1	0	1	0
388	1	0	1	0
389	69	50	63	50
390	1	0	1	0
391	3	83	3	83
392	2	9	2	9
393	1	0	1	0
394	1	0	1	0
395	16	3	270	16
396	1	0	1	0
397	2	0	2	0
398	1	0	154	7
399	1	0	1	0
400	9	- 10	9	- 10
total		374		454

7 real- AdaBoost i- AdaBoost LF1

	T=30	T=50	T=70	T=100
351	18	12	13	13
352	63	63	63	63
353	20	20	20	20
354	13	14	13	13
355	0	0	0	0
356	1	1	1	1
357	34	34	34	34
358	0	0	0	0
359	0	0	0	0
360	1	1	1	1
361	0	0	0	0
362	0	0	0	0
363	0	0	0	0
364	0	0	0	0
365	1	0	0	0
366	3	3	3	3
367	6	6	6	6
368	0	0	0	0
369	- 2	0	- 2	- 2
370	0	0	0	0
371	0	0	0	0
372	6	6	6	6
373	0	0	0	0
374	33	33	24	33
375	0	0	0	0
376	0	0	0	0
377	0	0	0	0
378	0	0	0	0
379	0	0	0	0
380	0	0	0	0
381	0	0	0	0
382	0	0	0	0
383	87	84	84	84
384	0	0	0	0
385	12	22	5	11
386	0	0	0	0
387	0	0	0	0



	T=30	T=50	T=70	T=100
388	0	0	0	0
389	44	50	53	52
390	10	0	0	0
391	83	83	83	83
392	9	9	9	9
393	0	0	0	0
394	0	0	0	0
395	13	16	17	10
396	0	0	0	0
397	0	0	0	0
398	7	7	9	9
399	0	0	3	0
400	- 10	- 10	- 10	- 10
total	442	454	435	439

8

7 i- AdaBoost real- AdaBoost LF1  
 LF1  
 374 385 i- AdaBoost가  
 real- AdaBoost 가

가 0  
 가 8 i- AdaBoost  
 i- AdaBoost  
 LF1 30, 50, 70, 100  
 50 가 real- AdaBoost

9 TREC- 8 가 7  
 LF1 LF1 , i- AdaBoost (T=50) LF1  
 LF1 i- AdaBoost (T=50)가 가

	i- AdaBoost	PLt8f1	PLt8f2	CL99bf1l	pirc9BF1	Mer8BaLF1	AntBatch1	Scal8Ft	Maximal
351	12	22	25	20	25	- 320	15	0	51
352	63	82	115	12	130	100	9	149	570
353	20	32	42	4	28	21	- 1001	0	87
354	14	- 9	- 9	- 2	- 7	- 15	- 261	0	147
355	0	0	0	0	0	0	- 324	0	3
356	1	7	7	- 5	4	- 345	0	- 2	45
357	34	47	41	34	59	46	- 43	3	183
358	0	0	0	0	0	0	0	0	6
359	0	- 5	- 5	0	0	0	- 992	- 8	90
360	1	0	0	0	0	3	- 1002	0	36
361	0	0	0	- 6	0	0	0	0	6
362	0	0	0	1	0	0	0	- 2	15
363	0	0	0	- 2	- 4	- 4	0	0	12
364	0	0	0	0	0	- 2	- 496	0	6
365	0	- 16	- 28	- 2	4	0	3	0	9
366	3	0	0	4	3	22	- 1002	0	51
367	6	0	5	0	4	4	- 28	0	90
368	0	0	0	0	0	0	- 55	0	15
369	0	0	0	0	0	0	0	0	0
370	0	- 20	- 20	- 2	- 13	- 4	- 389	0	93
371	0	0	0	0	0	0	0	0	3
372	6	3	3	- 6	1	0	0	- 4	36
373	0	0	0	13	0	0	0	0	36
374	33	5	5	- 4	- 19	12	- 152	- 2	177
375	0	- 2	1	1	3	3	3	0	30
376	0	3	3	- 9	0	0	- 366	0	51
377	0	6	6	9	2	3	0	0	30
378	0	- 25	- 14	- 12	- 12	- 11	- 11	5	132
379	0	0	0	- 8	0	0	0	0	0
380	0	0	0	0	0	0	0	0	3

	i- AdaBoost	PLt8f1	PLt8f2	CL99bfl1	pirc9BF1	Mer8BaLF1	AntBatch1	Scai8Ft	Maximal
381	0	0	0	- 6	0	0	0	0	6
382	0	- 2	- 2	7	- 2	0	- 1002	0	12
383	84	97	113	- 2	72	27	- 7	- 2	201
384	0	0	0	0	0	0	0	0	9
385	22	9	14	24	28	19	11	0	81
386	0	0	0	0	- 2	0	0	0	18
387	0	- 10	- 10	- 6	- 1	3	- 1002	0	27
388	0	0	0	- 3	0	0	- 323	0	42
389	50	145	113	122	76	31	176	218	387
390	0	- 4	- 4	0	0	3	- 2	- 4	204
391	83	36	- 18	57	- 35	- 16	- 22	- 138	381
392	9	- 32	- 32	- 5	2	- 5	- 21	4	96
393	0	7	7	- 2	1	3	- 14	0	15
394	0	0	0	- 8	0	0	0	0	15
395	16	- 8	- 1	- 6	- 50	- 1	- 115	- 9	186
396	0	1	1	- 6	4	0	0	0	15
397	0	- 2	- 2	0	0	- 6	0	0	15
398	7	9	9	- 1	0	4	- 3	0	30
399	0	0	0	4	3	0	0	0	27
400	- 10	- 13	- 8	3	- 9	5	- 20	0	48
total	454	363	357	212	295	- 420	- 8436	208	3828

9 TREC- 8 가 i- AdaBoost LF1

i- AdaBoost가 LF1  
 i- AdaBoost  
 가 가 가  
 . 가 i- AdaBoost가 가  
 real - AdaBoost  
 .  
 가 .  
 352  
 i- AdaBoost가

	AdaBoost- tf	real- AdaBoost	i- AdaBoost
	III 500mhz, Windows98		
	100	100	30*5, 50*5, 70*5, 100*5
	5	5	7 , 9 , 11 , 15

# 6

TREC  
, Freund Schapire

가 가  
TREC 가 가

Reuter

가

가

10 가

TREC

가

TREC

LF1

가

가

가

1.

가

가

가

가

가

가

가

2. TREC 가

가

TF-IDF

3,000

가

3. i-AdaBoost

,

가

가

가  
가

가

가

가

,

가

[abney99] Steven Abney, Robert E. Schapire, and Yoram Singer. Boosting applied to Tagging and PP attachment. *Joint SIGDAT Conference on Empirical Methods in Natural Language Processing and Very Large Corpora*, 1999.

[bauer99] Bauer E., and Kohavi R.. An empirical comparison of voting classification algorithms: bagging, boosting, and variants. *Machine Learning*, volume 36 Issue 1, pages 105–139, 1999.

[breiman96a] L. Breiman. Bagging predictors. *Machine Learning*, volume 24, 123–140, 1996.

[breiman96b] L. Breiman. Bias, variance, and arcing classifiers. Technical Report 460. Berkeley, CA: University of California: Department of Statistics, 1996.

[breiman97] L. Breiman. Arcing the edge. Technical Report 486. Berkeley, CA: University of California: Department of Statistics. 1996.

[freund96a] Y. Freund, R. E. Schapire. Experiments with a new boosting algorithm. *Proceedings of the Thirteenth International Conference on Machine Learning*, pages 148–156, 1996.

[freund96b] Y. Freund, R. E. Schapire. Game theory, on-line prediction and boosting, *Proc. 9th Annu. Conf. on Comput. Learning Theory*, pp. 325-332, ACM Press, New York, NY, 1996.

[freund97] Yoav Freund, Robert E. Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of Computer and System Sciences*, volume 55, pages 119–139, 1997.

[freund99] Yoav Freund, Robert. E. Schapire. The large margin classification using the perceptron algorithm. *Machine Learning*, volume 37, Issue 3, pages 277–296, 1999.

[friedman] J. Friedman, T. Hastie, and R. Tibshirani. Additive logistic regression: A statistical view of boosting. *Annals of Statistics*, to appear.

[hull96] David Hull, Jan Pedersen, and Hinrich Schutze. Method combination for document filtering. *Proceeding of the Nineteenth Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 279–288, 1996.

[lewis92] David Lewis. An evaluation of phrasal and clustered representations on a text categorization task. *Proceedings of the Fifteenth Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 37–50, 1992.

[lewis94] David Lewis and William Gale. A sequential algorithm for training text classifiers. *Proceedings of the Seventeenth Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 3–12, 1994.

[lewis96] David Lewis, Robert Schapire, James Callan, and Ron Papka. Training algorithm for linear text classifiers. *Proceeding of the Nineteenth Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 298–306, 1996.

[neuralnet] Simon Haykin, *Neural Networks*. Prentice-Hall, Inc, 1999.

[quinlan96] J. R. Quinlan. Bagging, Boosting, and C4.5. *Proceedings of the Thirteenth National Conference on Artificial Intelligence*, pages 725–730, 1996.

[schapire98a] R. E. Schapire, Y. Freund, P. Bartlett, W. S. Lee. Boosting the margin: A new explanation for the effectiveness of voting methods. *The Annals of Statistics*, volume 26, 1651–1686, 1998.

[schapire98b] Robert E. Schapire, Yoram Singer, and Amit Singhal. Boosting and Rocchio applied to text filtering. *Proceedings of the 21st Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, Query and Profile Modification, pp. 215-223, 1998.

[schapire99a] Robert E. Schapire, Yoram Singer. Improved boosting algorithms using confidence-rated predictions. *Machine Learning*, volume 37, Issue 3, 1999.



[schapire99b] Robert E. Schapire, and Yoram Singer. BoosTexter: a boosting-based system for text categorization. *Machine Learning*, 1999.

[webb] Geoffrey I. Webb. MultiBoosting: A technique for Combining Boosting and Wagging. *Machine Learning*, to appear.

[yang94] Yiming Yang. Expert network: effective and efficient learning from human decisions in text categorization and retrieval. *Proceeding of the Seventeenth Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 13–22, 1994.

[yang95] Yiming Yang. Noise reduction in a statistical approach to text categorization. *Proceeding of the Eighteenth Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 256–263, 1995.

## Abstract

Text filtering is a task of deciding whether a document has relevance to a specified topic. As Internet and Web becomes wide-spread and the number of documents delivered by e-mail explosively grows the importance of text filtering increases as well. The aim of this thesis is to improve the accuracy of text filtering systems by using machine learning techniques. We apply AdaBoost algorithms to the filtering task. An AdaBoost algorithm generates and combines a series of simple hypotheses. Each of the hypotheses decides the relevance of a document to a topic on the basis of whether or not the document includes a certain word. We begin with an existing AdaBoost algorithm which uses weak hypotheses with their output of 1 or -1. Then we extend the algorithm to use weak hypotheses with real-valued outputs to improve error reduction rate and final filtering performance. Next, we attempt to achieve further improvement in the AdaBoost's performance by first setting weights randomly according to the continuous Poisson distribution, executing AdaBoost, repeating these steps several times, and then combining all the hypotheses learned. This has the effect of mitigating the overfitting problem which may occur when learning from a small number of data. Experiments have been performed on the real document collections used in TREC-8, a well-established text retrieval contest. This dataset includes Financial Times articles from 1992 to 1994. The experimental results show that AdaBoost with real-valued hypotheses outperforms that with binary-valued hypotheses, and that AdaBoost iterated with random weights further improves filtering accuracy. Comparison results of all the participants of the TREC-8 filtering task are also provided.

